



## Bias In, Bias Out?

How to Fight Biases in AI and Create Trust

**Dr. Gerald Fahner**

Analytic Science-Sr Principal Scientist, FICO

Algorithmic prejudice

# Facebook's ad system seems to discriminate by race and gender

*New research shows that Facebook's ad-distribution software is disturbingly biased*

THE DAILY NEWSLETTER

Sign up to our daily email newsletter

# NewScientist

News Technology Space Physics Health Environment Mind Video | Travel Events Jobs

## Face-recognition software is perfect – if you're a white man



TECHNOLOGY 13 February 2018



Business

Markets

World

Politics

TV

More

BUSINESS NEWS

OCTOBER 9, 2018 / 10:12 PM / 6 MONTHS AGO

# Amazon scraps secret AI recruiting tool that showed bias against women

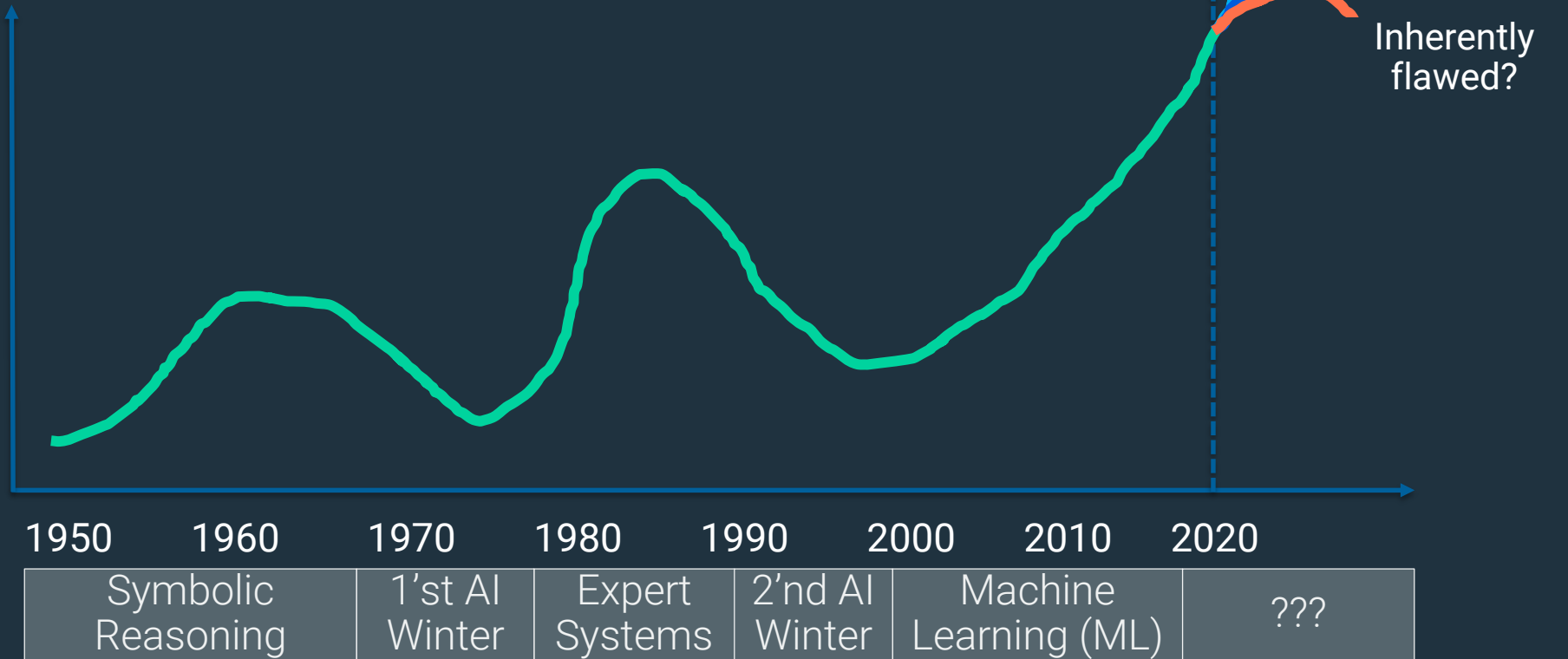
Jeffrey Dastin

8 MIN READ



# Up's and Down's of AI – What's Next to Come?

Impact  
Pervasiveness



# Gartner Hype Cycle for Emerging Technologies, 2019

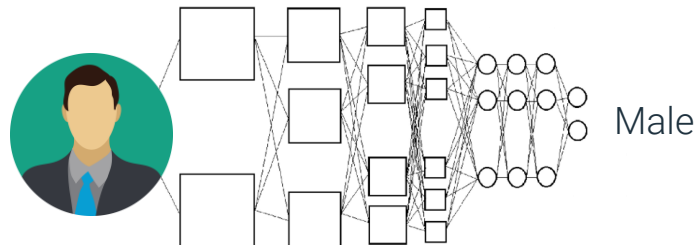


# Bias Culprit 1: Using Training Data That Lack Diversity

Training Data Labelled by Gender



Deep NN for gender classification



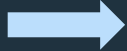
Error rates  
for various commercial products\*

Dark Females	24% to 36%
Light Females	0% to 2%
Dark Males	0% to 6%
Light Males	0.0% to 0.3%

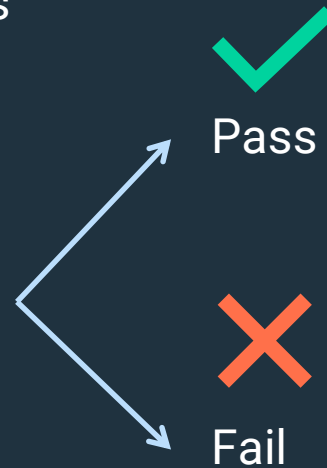
\*J. Buolamwini and T. Gebru, Gender Shades: "Intersectional Accuracy Disparities in Commercial Gender Classification", Proceedings of Machine Learning Research 81:1–15, 2018

## Bias Culprit 2: Training a Machine With Subjective Labels

Resumes



Judgmental prescreen decisions





## Bias Culprit 3: Letting Irrelevant “Big Data” Features Creep Into a Model

Resumes



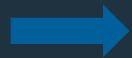
Zillion features

	innovate	softball	influence	women	+ many more
Resume 1	1	0	0	1	...
Resume 2	0	0	1	2	...
Resume 3	3	1	0	0	...

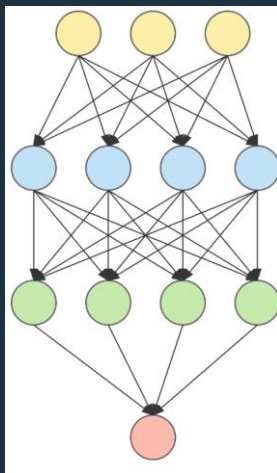
## Bias Culprit 4: Prioritizing “Best Fit” Over Transparency



Train model



Zillion  
features

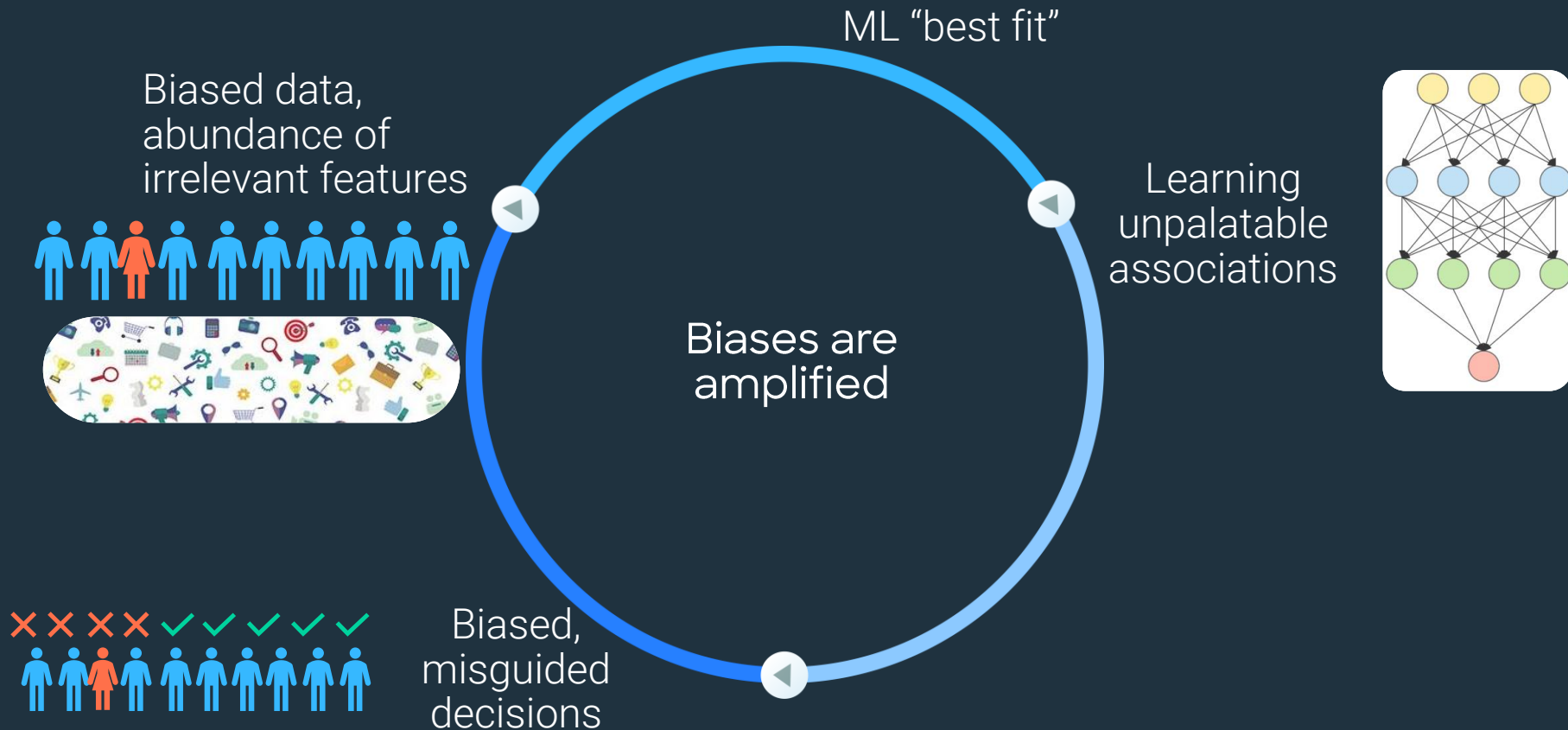


Pass/Fail

Biased

Feature	Impact on 'Pass' decision
innovate	+1.37
softball	-0.22
influence	+0.28
women	-0.31
execute	+0.57
...	...

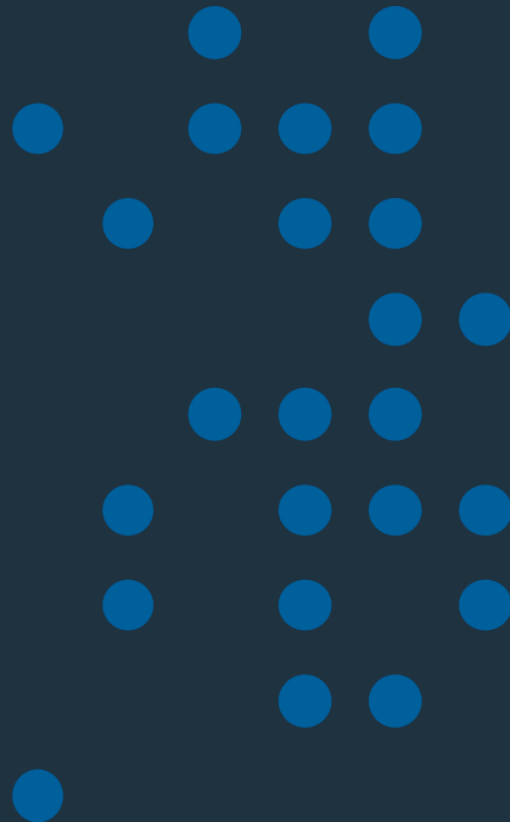
# Vicious Artificial Intelligence Cycle



## Questions for You

How can we mitigate algorithmic bias (discrimination) against demographic groups?

- A. Use all available data sources to build the most accurate model
- B. Don't use demographic group membership as a predictor



In the US, if a credit scoring model:

Doesn't use prohibited factors or proxies  
for prohibited factors

AND

Is empirically derived, demonstrably and  
statistically sound ("EDDSS")

THEN

The model affords certain protections  
against disparate treatment and disparate  
impact claims



## Questions for You

How can we mitigate algorithmic bias (discrimination) against demographic groups?

- A. Use all available data sources to build the most accurate model
- B. Don't use demographic group membership as a predictor
- C. Allow score cutoffs to differ between demographic groups, such that equal proportions from each group are selected
- D. Collect relevant and reliable variables that are predictive for all groups

Collect relevant and reliable variables  
that are predictive for all groups

# Bias Buster 1: Prioritize Plausibly Causal Variables Predictive for All Groups

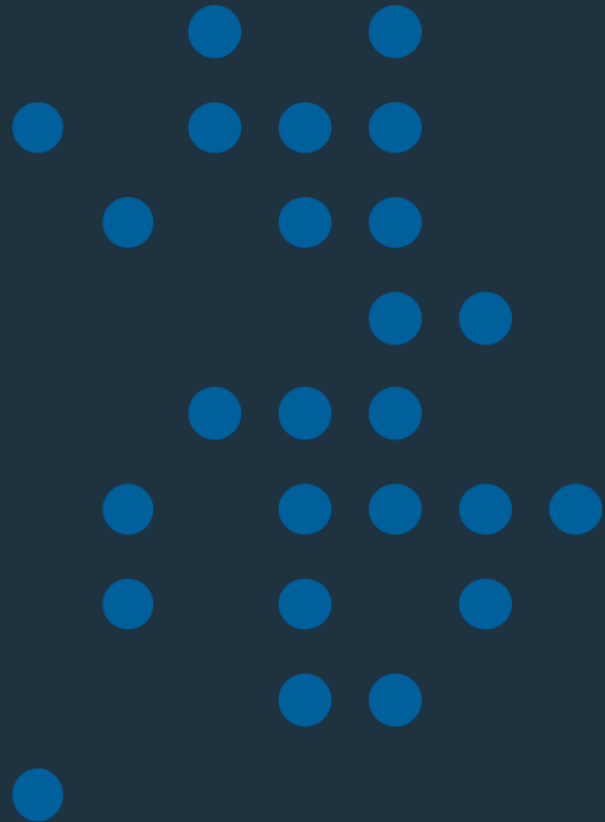




# Questions for You

Do you believe the following situations are fair?

- A. A parole algorithm is equally predictive for black and white defendants
- B. A parole algorithm falsely flags black defendants twice as likely as future criminals, than white defendants
- C. A university admits 44% of male applicants but only 35% of female applicants



# Berkeley Admissions Decisions\*

	Men		Women	
	Applicants	Admitted	Applicants	Admitted
Total	8442	44%	4321	35%

Biased against women?

Department	Men		Women	
	Applicants	Admitted	Applicants	Admitted
A	825	62%	108	82%
B	560	63%	25	68%
C	325	37%	593	34%
D	417	33%	375	35%
E	191	28%	393	24%
F	373	6%	341	7%

Not biased on department level

Women applied more to departments that are hard to get into

\*P. Bickel, E. Hammel and J. O'Connell: "Sex Bias in Graduate Admissions: Data From Berkeley", Science. 187 (4175): 398–404, 1975

# ~~Vicious Artificial~~ Enlightened Augmented Intelligence Cycle

ML combined with human insight

- Representative
  - Reliable
  - Plausibly causal
- Well-designed data

More palatable models

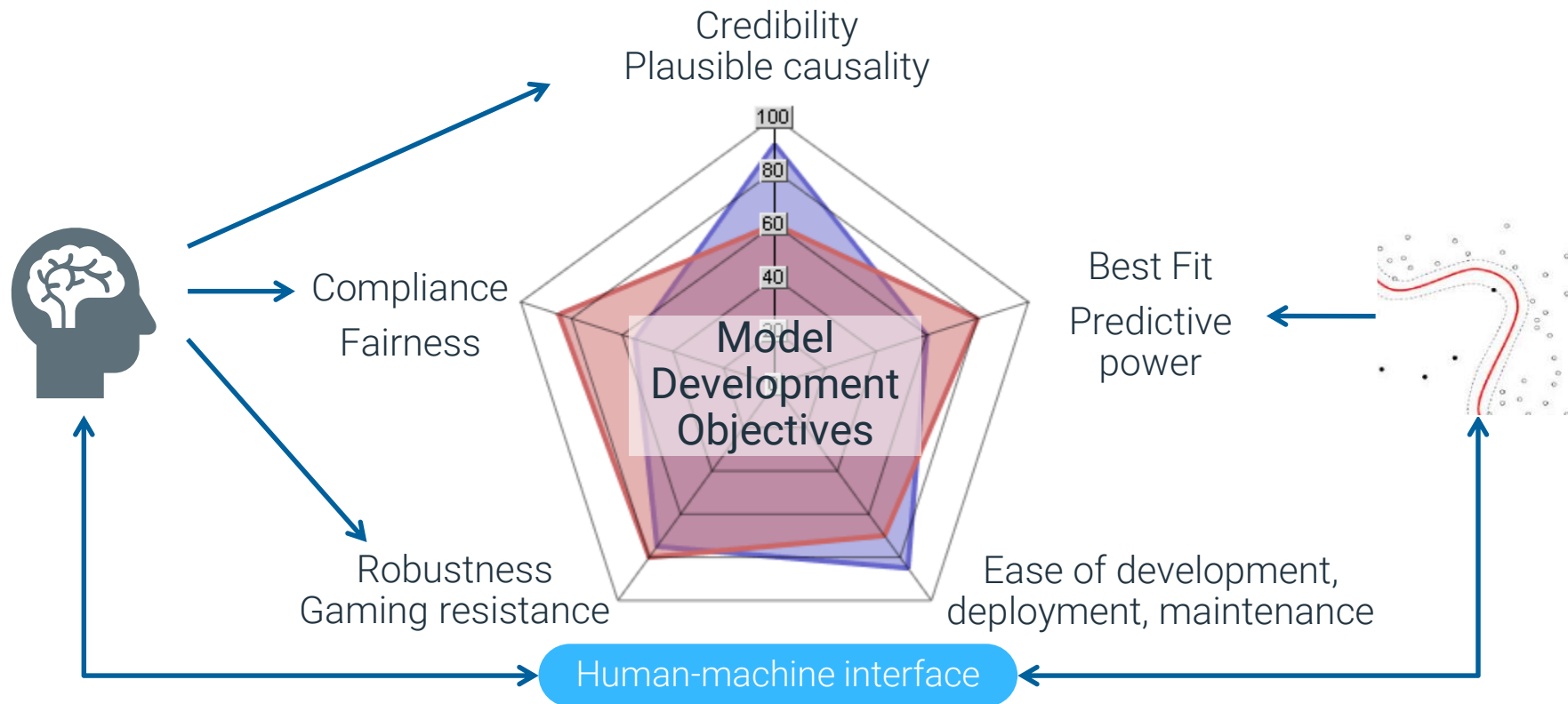
- Compliant
- Credible
- Plausibly causal

Biases are reduced

- Compliant
  - Actionable
  - Fair
- Less biased, reasoned decisions, feedback to people



## Bias Buster 2: Combine Machine Learning With Human Insight

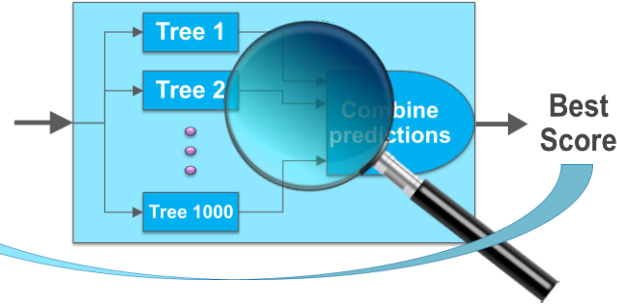
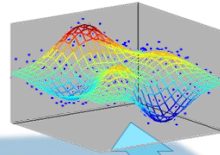


# Human-Machine Interface for Highly Effective Scorecard Development

[https://www.thinkmind.org/index.php?view=article&articleid=data\\_analytics\\_2018\\_1\\_30\\_60077](https://www.thinkmind.org/index.php?view=article&articleid=data_analytics_2018_1_30_60077) (as of 10/17/2019)

- 1 Develop “best fit” black box AI/ML model
- 2 Seek to understand black box
- 3 Augment data with **Best Score**
- 4 Approximate **Best Score** by automatically grown, explainable scorecard system

## Training Data



## Diagnostics:

- Predictive performance benchmark
- Variable importance
- Nonlinear patterns, interactions

- 5 Hone final scorecards to warrant compliance

## Full Population

### Delinquent

### Current

Util < 40% .0 Util >= 40%

Number of months since the most recent serious delinquency	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55
Overall utilization on revolving trades	No revolving trades Under 6% 7 - 19% 20 - 49% 50 - 89% 90% or more	30 65 45 25 15
Number of months since the most recent serious	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55

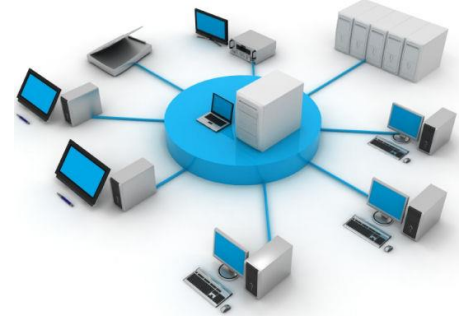
Number of inquiries in last 6 mos.	0 1 2 3 4+	70 60 45 25 20
Number of bankrupt trade lines	0 1 2 3 4+	15 25 55 60 50
		15 15 25 60 50

Number of months since the most recent serious delinquency	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55
Overall utilization on revolving trades	No revolving trades Under 6% 7 - 19% 20 - 49% 50 - 89% 90% or more	30 65 45 25 15
Number of months since the most recent serious	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55

Number of months since the most recent serious delinquency	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55
Overall utilization on revolving trades	No revolving trades Under 6% 7 - 19% 20 - 49% 50 - 89% 90% or more	30 65 45 25 15
Number of months since the most recent serious delinquency	No serious delinquency 0 - 5 6 - 11 12 - 23 24+	75 10 15 25 55
Overall utilization on revolving trades	No revolving trades Under 6% 7 - 19% 20 - 49% 50 - 89% 90% or more	30 65 45 25 15
Number of months in file	Below 12 12 - 23 24 - 47 48 or more	12 35 60 75

**APPROVED**

- 6 Deploy multi-scorecard system



# Human-Machine Interface Dominates Other Approaches

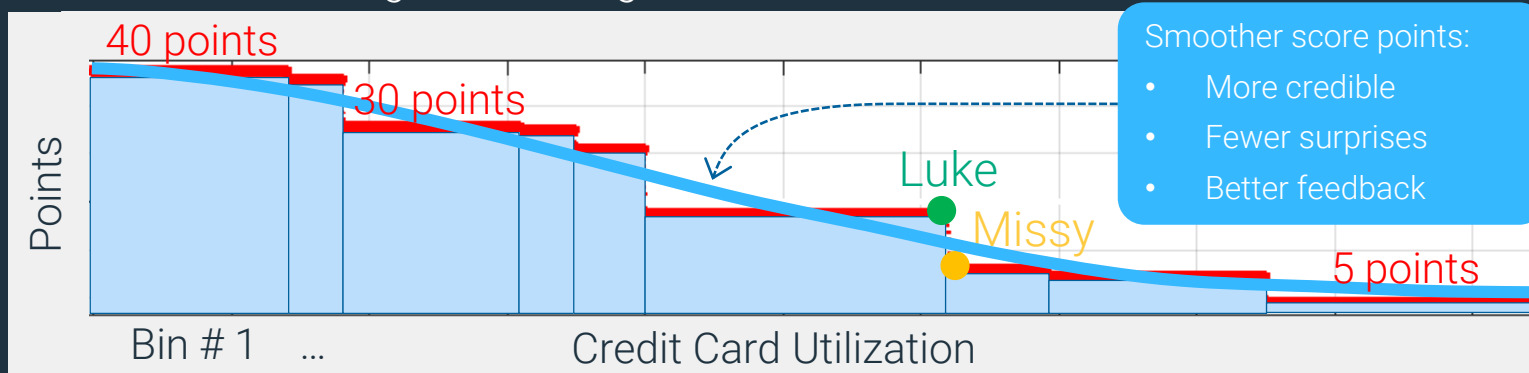
Model / Technology*	Development Effort (Resource Hours)	Explainability / Credibility	Predictive power improvement** over FICO® Score 9
FICO® Score 9	Laborious	High	N/A
Gradient Boosting	Modest	Limited	+1.7%
Multilayer Neural Net	Modest	Limited	+0.5%
Human-Machine Interface	Modest	High	+0.3%

\*Same data used to train and test all models

\*\*Performance metric: KS (Kolmogorov-Smirnov statistic)

# Further Improving Explainability With “Liquid” Scorecards

## Binning and Scoring of a Continuous Behavior Variable



	Credit Card Utilization	Score points
Luke	60%	18 Points
Missy	61%	8 Points

*“Because the algorithm said there’s a bin break here” ?*

# Prudently Applied AI/ML for Credit Scoring Brings Benefits at Scale

1950

Subjective decisions / Rampant discrimination

Credit scoring eliminates judgment

Regulation protects vulnerable groups

FICO® Score for diverse populations

2019

Greater financial inclusion with alternative data (UltraFICO™ Score, FICO® Score XD)

We keep innovating to reduce bias and boost fairness



“Biased AI” is not an inherent flaw.

Human and artificial intelligence can  
combine to fight biases at a massive scale.

If it sounds easy, it isn't.  
But it's possible, and it's necessary.



FICO®

Thank You!

Dr. Gerald Fahner